**WHAT IS CLAIMED IS:**

1. A service distribution device for distributing services among a plurality of

servers on a network to balance the server loads, comprising:

5      a packet capture device capturing packets transmitted through the network;

a server identifier recording information pertaining to the captured packets into a server

log for each server;

a service identifier recording information pertaining to the captured packets into a

service log for each service;

10      a server modeling module setting up the simulation model for each server from the

server records;

a service modeling module setting up the simulation model for each service from the

service records;

a simulator reading in the server model and the service model and running each

15  simulation; and

a server selection module selecting and specifying an optimum server to distribute

services to based on a simulator result.


2.      The service distribution device of claim 1, further comprising a packet relay

20  device obtaining packets using a packet capture module mounted on said packet relay device,

which relays packets between a client and the servers.

3.　　The service distribution device of claim 1,

wherein said server modeling module constructs a server model having a queue

corresponding to a transmission process using the server log and a server transmission

throughput, a server processing time, and a unit processing time as parameters,

5　　　　wherein the server transmission throughput is calculated from a total size L of

an arbitrary, continuous string of the continuously transmitted packets using the

formula $L / (t_e - t_s)$ where $t_e$ is an ending packet capture time and $t_s$ is a starting

packet capture time, and

wherein the server processing time is calculated using the formula

10　$(t_s - t_c) - (l_s + l_c) / B$, wherein $t_s$ and $l_s$ are the capture time and size of a server

response packet, respectively, $t_c$ and $l_c$ are the capture time and size of a corresponding

client response packet, respectively, and B is a network speed.

4.　　The server distribution device of claim 1, wherein said service modeling module

15　calculates the following parameters from the service log by constructing a service model for

each service:

a ratio of the number of sessions for each service to the number of sessions for all

services,

a session starting frequency or time interval,

20　　　a number of transmissions between the client and server per session,

a client response size, packet size, and packet count per transmission,

a server response size, packet size, and packet count per transmission, and

a time from the server response until the client response.

5.    The service distribution device of claim 1, wherein said simulator performs a simulation using the server model and the service model and generates a mean value or a median value of a session time for the specific service.

5

6.    The service distribution device of claim 1, wherein said server selection module determines a standard value using an output of a single simulation run for each service by said simulator, and determines that a high load state exists if a difference between, or the ratio of, the standard value and the output of the simulation of a plurality of sessions exceeds a pre-

10    determined threshold.

7.    The service distribution device of claim 6, wherein when said server selection module receives a server distribution query, said server selection module sets a server permission to be a starting frequency of the session that will cause a high load state for the

15    service in question for each server, and specifies a server having the biggest difference between the session starting frequency and the permission as a server for distribution.

8.    The service distribution device of claim 6, wherein when said server selection module receives a distribution server query, said server selection module runs a simulation for

20    a service in question for each server and specifies a server for which a result of a ratio for which $\beta$ multiplied by the standard value is less than or equal to $\gamma$.

9.    The service distribution device of claim 6, wherein when said server selection module receives a distribution server query, said server selection module runs a simulation for a service in question for each server and specifies as a distribution server, a server for which a result of ratio for which $\beta$ multiplied by the standard value is smallest.

10.    The service distribution device of claim 4, wherein said service modeling module categories each session transmission as a connection request and response, and a command transmission, a data transmission, a response, and an end, and calculates the parameters for each session transmission based upon category.

11.    The service distribution device of claim 7, wherein the permissions of each of the servers are taken as weighted values of a service distribution or relative ratios of the permissions are taken as server distribution ratios.

12.    A service distribution device for distributing services among a plurality of servers to balance server loads, comprising:

a server modeling module generating a simulation model for each server and a service modeling module generating a simulation model for each service;

a simulator reading the server models and the service models and running a plurality of simulations; and

a server selection module determining which servers have low loads based on results of the simulations and selecting the servers with low loads to receive the services.

13.    A method for distributing services among a plurality of servers to balance server

loads, comprising:

    generating a simulation model for each server and each service;

    running a plurality of simulations using the server and service models; and

5    determining which servers have low loads based on results of the simulations

and selecting the servers with low loads to receive the services.


14.    A computer-readable storage controlling a computer and comprising a

process of:

10    generating a simulation model for each server and each service;

    running a plurality of simulations using the server and service models; and

    determining which servers have low loads based on results of the simulations

and selecting the servers with low loads to receive the services.

15